# *All Tools Data Analysis Project*

During our training at

***National Telecommunication Insitute And ITIDA***

Full Data Analysis Process with Python

Full Data Analysis Process with SQL

Full Data Analysis Process with Excel

# *Our Dataset Talk About*

**Project Goal:** Analyze song performance, compare songs by
genre, language,
duration, and explicit content, and extract valuable insights to
support decision
making in the music industry

*Key dimensions covered in the dataset include:*

- **Song Attributes:** [song_title], [artist], [album], [genre],
          [language], [duration]

- **Performance Metrics:** [popularity], [stream]

- **Classifications:** [explicit_content], [popularity_level],
     [streams_level], [duration_minute], [date_group]

- **Production Data:** [composer], [producer]

# *Our Full Data Analysis Process with Excel*

## . Data cleaning

No duplication

    o Deal with outliers

        ▪ There is outliers in [duration] and its count 349 it will be

deleted

     o Deal with nulls

        ▪ Filling nulls in [language] by mode

        ▪ Filling nulls in [duration] by mean

        ▪ Delete [collaboration] Because it contains many nulls = 35000

o Feature Engineering

        ▪ Create columns like

            o popularity_level

             o duration_minute

              o streams_level

              o date_group

▪ the columns that created will Facilitate analysis and comparison

across different group

1-

# *Our Full Data Analysis Process with Excel*

**We also performed a complete data analysis process using Excel, leveraging pivot tables, slicers, and dashboards to ensure data cleaning, interactive exploration, and clear visualization of insights from the Spotify_songs data set**

2-**Pivot Tables & Slicers**
**Pivot Tables Built For:**
    1) Top Music Labels by Number of Songs
     2) Average Song Popularity by Label
     3) Average Streams: Old vs Recent Songs
     4) Average Popularity: Explicit vs Non-Explicit Songs

**Slicers Added:**

# *Our Full Data Analysis Process with Excel*

• Key Insights

• **Sony Music** and **Universal Music** are the top labels in terms of number of songs.

• **English songs** dominate the dataset, representing about **71%** of all songs.

• **Older songs** achieved the highest number and average of **streams** compared to newer releases.

• **Recent songs** (new/medium) attract fewer streams than older ones.

• **Sony Music** and **Universal Music** have the highest average popularity, while **Indie labels** show lower popularity.

• There is **no significant difference** in popularity between **Explicit** and **Non-Explicit** songs.

• **English songs** are the most common language for **Explicit content**.

• The largest number of releases occurred in the **old period (17,016 songs)**, while recent years show fewer releases (**16,147 songs**).

# Spotify_songs dashboard

**label**

Def Jam

Indie

Sony Music

Universal Music

Warner Music
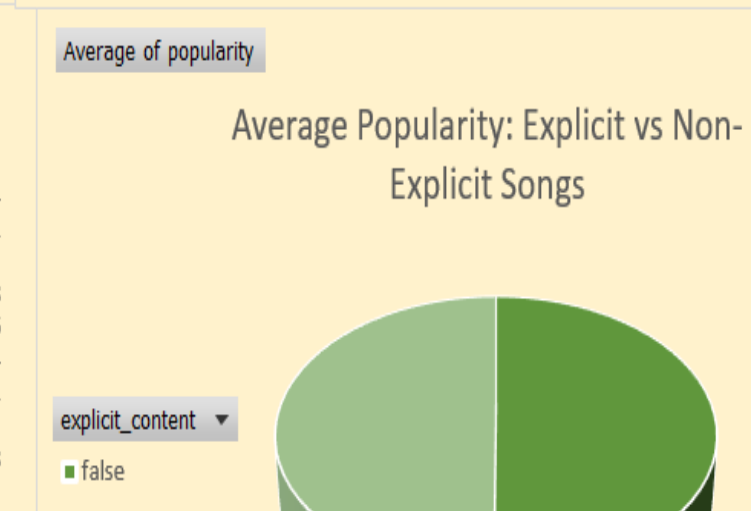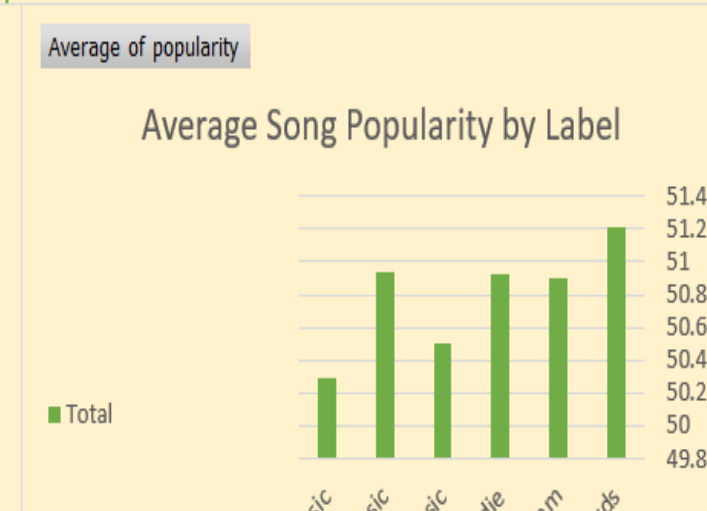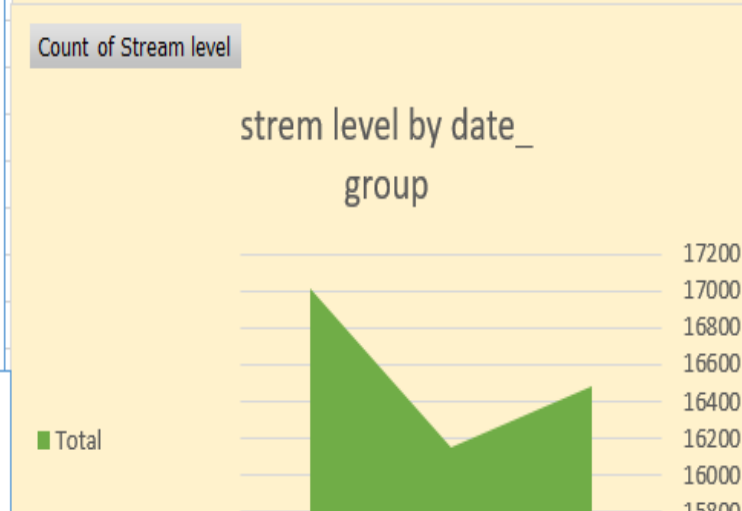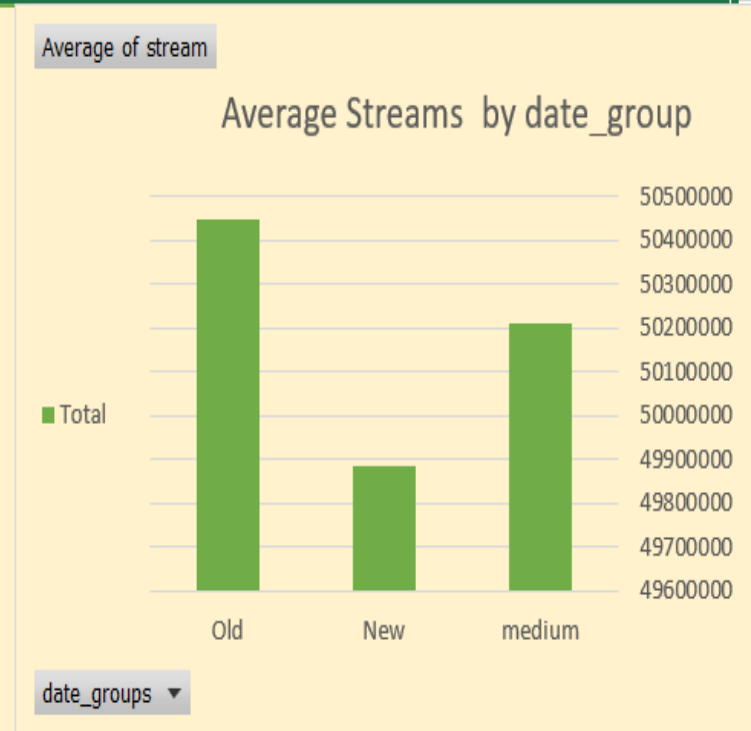
**language**

English

French

German

Italian

Japanese

Korean

Spanish

**stream**

1899

Count of song_id

## Top Music Labels by Number of Songs

Warner Music
Universal Music
Sony Music
label
Indie
Def Jam
Atlantic Records

■ Total

label ▼

8600 8400 8200 8000 7800

count of song_id

Count of song_id

## COUNT SONG BY LANGUGE

language ▼

■ English
■ French
■ German
■ Italian
■ Japanese
■ Korean
■ Spanish

9%
5%
3%
4%
3%
5%
71%

Average of stream

## Average Streams by date_group

50500000
50400000
50300000
50200000
50100000
50000000
49900000
49800000
49700000
49600000

■ Total

Old        New        medium

date_groups ▼

Count of Stream level

## strem level by date_ group

17200
17000
16800
16600
16400
16200
16000
15800

■ Total

Average of popularity

## Average Song Popularity by Label

51.4
51.2
51
50.8
50.6
50.4
50.2
50
49.8

■ Total

Average of popularity

## Average Popularity: Explicit vs Non-Explicit Songs

explicit_content ▼

■ false

# *Our Full Data Analysis Process with SQL*

We also performed a complete data analysis process using
*SQL* to ensure robust validation, efficient querying,
and powerful insights from the Spotify_songs dataset.

## DATA EXPLORATION

□ **Univariate Analysis:**
Studied each column separately. For numeric features
(`popularity,`
`stream`), we used **histograms** and **boxplots** to check
distributions and
outliers. For categorical features (`genre, language,`
`duration_minute, explicit_content, date_group,`
`popularity_level, streams_level`),

# *Our Full Data Analysis Process with SQL*

**Data cleaning**

No duplication

 o Deal with outliers

   ▪ There is outliers in [duration] and its count 349 it will be

deleted

  o Deal with nulls

   ▪ Filling nulls in [language] by mode

   ▪ Filling nulls in [duration] by mean

   ▪ Delete [collaboration] Because it contains many nulls = 35000

o Feature Engineering

  ▪ Create columns like

    o popularity_level

    o duration_minute

    o streams_level

    o date_group

▪ the columns that created will Facilitate analysis and comparison

across different group

# *Our Full Data Analysis Process with SQL*

Feature Engineering •

▪ Create columns like •

o popularity_level •

o duration_minute •

o streams_level •

o date_group •

Our Full Data Analysis Process with SQL

# *Our Full Data Analysis Process with SQL*

•Key Insights

•**English songs** dominate the dataset, representing about **71%** of all songs.
•**Older songs** achieved the highest number and average of **streams** compared to newer releases.
•**Recent songs** (new/medium) attract fewer streams than older ones.
•**Sony Music** and **Universal Music** have the highest average popularity, while **Indie labels** show lower popularity.
•There is **no significant difference** in popularity between **Explicit** and **Non-Explicit** songs.
•**English songs** are the most common language for **Explicit content**.
•The largest number of releases occurred in the **old period (17,016 songs)**, while recent years show fewer releases (**16,147 songs**).

# *Our Full Data Analysis Process with Python*

Import Libraries & Load Files

Imported necessary libraries (Pandas, NumPy, Matplotlib, etc.)

Loaded dataset files into DataFrame

Data Exploration:

  - Checked dataset shape (number of rows and columns)

  - Checked column data types and null values

   - Viewed last rows to inspect data

# Data Cleaning

Null Values: •

- Checked number of null values in each column

- Dropped rows where all values were null

- Dropped the 'collaboration' column

- Replaced remaining null values with mode (categorical) or mean (numerical)

Duplicates:

- Checked for duplicate rows

- Found no duplicates

Group By & Pivot Table:

- Used groupby() and pivot_table() to summarize and aggregate data

# Insights / Key Questions

Q1: Which songs are the most popular?

Q2: Which songs have the highest number of streams?

Q3: Which genre has the highest average popularity?

Q4: Which genre has the highest average streams?

Q5: Are explicit songs more popular than non-explicit songs?

Q6: Do major labels (Universal, Sony, Warner, Def Jam) have higher streams?

Q9: Which languages are most used in songs?

Q10: Which languages have the highest average popularity?

Q11: Which artists released the highest number of songs?

Q12: Which artists achieved the highest total streams?

Q19: Which producers or composers are associated with the most popular songs?

Q20: Which genres produce more explicit songs? •

# Data Visualization

Bar Charts – Important Categorical Columns: Shows distribution of genres, explicit content, languages, albums, and artists.

Histograms – Numerical Columns: Displays distribution of duration, popularity, and streams.

Pie Chart – Language Distribution: Shows proportion of songs in different languages.